# An integrative approach to investigate natural variation in the accumulation of aliphatic glucosinolates in *Arabidopsis thaliana*

Suraj Sharma[1,2], Ovidiu Popa[2], Stanislav Kopriva[1,3], Oliver Ebenhoeh[1,2]

[1] Cluster of Excellence on Plant Sciences (CEPLAS), Heinrich-Heine University Duesseldorf, Germany

[2] Institute for Theoretical Biology, Heinrich-Heine University Duesseldorf, Germany

[3] Botanical Institute, University of Cologne, Germany

## Abstract

Glucosinolates are a fascinating class of specialised metabolites found in the plants of *Brassicacea* family. The variation in glucosinolate composition across different Arabidopsis ecotypes could be a result of allelic compositions at different biosynthetic loci. The contribution of methylthioalkylmalate synthase (MAM) genes to diversity of glucosinolate profiles across different Arabidopsis ecotypes has been confirmed by genetic analyses. Different MAM isoforms utilise different chain-elongated substrates for glucosinolate biosynthesis causing thus a variation in chain lengths across different Arabidopsis ecotypes. To further investigate the relationship between the genotype and the associated metabolic phenotype, we studied the diversity of genes and enzymes of glucosinolate biosynthesis. Using Shannon entropy as a measure we revealed that several genes of the pathway show a clear derivation from the expected behaviour, either accumulating non-synonymous SNPs or showing signs of purifying selection. We found that the genotype-phenotype relationship is much more complicated than inferred from the diversity of MAM synthases. We conclude therefore, that the ON/OFF feature of key QTLs is not enough to elucidate the diversity of glucosinolates across different *Arabidopsis thaliana* ecotypes and that glucosinolate profiles are determined also through the polymorphic residues along the coding regions of multiple metabolic genes.

## Introduction

Plant secondary metabolism produces a huge variation in structures and molecules with a plethora of functions for the plants but also for human nutrition and health (Owen, Patron, Huang, & Osbourn, 2017). One class of such secondary metabolites are glucosinolates in the Brassicaceae. Glucosinolates (GSLs) are important for the plants as a defence against herbivores, fungi, and other pathogens (Halkier & Gershenzon, 2006). They are also determinants of taste and flavour of cruciferous vegetables and responsible for their positive health properties (Traka & Mithen, 2009). More than 140 different GSL structures have been described, with a great variation not only between species, but also among ecotypes of the same species (Agerbirk et al., 2015; Clarke, 2010). GSLs are synthesised from amino acids and are divided into three classes: aliphatic GSLs, derived from alanine, leucine, isoleucine, valine and methionine (Met), indolic GSLs synthesised from tryptophan, and aromatic GSLs from phenylalanine (Phe) (Halkier & Gershenzon, 2006; Sønderby, Geu-Flores, & Halkier, 2010). All GSLs possess the same core structure, which comprises a glucose residue linked to a (Z)-N-hydroximic sulfate ester via a sulphur atom (Fahey, Zalcmann, & Talalay, 2001). GSL biosynthesis consists of three independent steps: (i) chain elongation of selected precursor amino acids (only Met and Phe), (ii) formation of the core GSL structure, and (iii) secondary modifications of the amino acid side chain. The diversity of GSLs derives from the side-chain elongation and secondary modifications.

GSLs are best known for their function in interaction between plants and herbivores. Upon tissue damage, the GSLs stored in the vacuoles come in contact with the enzyme myrosinase, which cleaves the thio-glucoside bond. The resulting aglycones are unstable and form volatile isothiocyanates or nitriles (Halkier & Gershenzon, 2006). Depending upon the herbivore, the volatiles of specific GSLs can act as feeding deterrents or stimulants (Buskov, Serra, Rosa, Sørensen, & Sørensen, 2002; Gabrys & Tjallingii, 2002; Lazzeri, Curto, Leoni, & Dallavalle, 2004; Mewis, Ulrich, & Schnitzler, 2002; Miles, Campo, & Renwick, 2005; Noret et al., 2005). Therefore, there is often an overlap between quantitative trait loci (QTL) for GSL accumulation and insect resistance (Kroymann, Donnerhacke, Schnabelrauch, & Mitchell-Olds, 2003). A possible outcome of this heterogeneous natural selection on GSLs is the quick evolution of new compounds or new patterns of compound accumulation (Daxenbichler et al., 1991; Rodman, 1980). New GSLs may increase resistance to herbivores that have become adapted to existing defences, whereas new patterns of GSL accumulation

60  may provide a unique complement of defences by slowing down the counter-adaptation of
61  herbivores.

62      The GSL defence system is one of the few systems where systematic data assessing
63  between and within species variation at both phenotypic and causal genetic level is available
64  (Halkier & Gershenzon, 2006; Sønderby et al., 2010). Natural variation within or between
65  species is regulated by a complex network of genes and associated polymorphisms (Fisher,
66  1930; Kliebenstein, Kroymann, et al., 2001; Lynch, Walsh, & others, 1998). These variations,
67  however, complicate our understanding of how certain genes behave in context of a species
68  as we often study a single genotype. Thus, understanding a metabolic pathway requires
69  studies involving more than one ecotype. For example, methionine derived aliphatic GSLs
70  differ in length of the side chain caused by the elongation of the amino acid, as well as by
71  further modifications, e.g. oxidation of the sulfur atom (Halkier & Gershenzon, 2006;
72  Sønderby et al., 2010). In the model plant *Arabidopsis thaliana*, aliphatic GSLs of six
73  different chain-lengths, referred to as 3C to 8C GSLs, but with different side chains are
74  found. The diversity in length of aliphatic GSLs can be explained by the variation in the
75  iterative chain-elongation cycles, each adding one methylene group to the Met and/or
76  elongated Met molecule (Halkier & Gershenzon, 2006). The QTL responsible for
77  determining the chain-elongation of GSLs is GS-ELONG (Magrath et al., 1994). GS-ELONG
78  is highly variable across different Arabidopsis ecotypes, with common large indel
79  polymorphism (Kroymann et al., 2003). The gene underlying the GS-ELONG QTL is
80  methylthioalkylmalate synthase (MAM3), encoding the key enzyme of the chain elongation
81  cycle (Kroymann et al., 2001). However, the GS-ELONG locus harbours two more genes,
82  isoforms of MAM3 called MAM1 and MAM2. MAM3 is present in all Arabidopsis ecotypes,
83  and some ecotypes possess both additional genes whereas others possess either MAM1 or
84  MAM2, with some, such as *Landsberg erecta* (*Ler*), having a truncated (non-functional)
85  MAM1 in addition to MAM2 (Kroymann et al., 2003). While the presence of functional
86  MAM genes has been described as responsible for the variation in aliphatic GSL side chain
87  length, very little is known about contribution of other genetic variation, mainly single
88  nucleotide polymorphism (SNP), at this locus.

89      In this paper, we investigate the link between the diversity of GSL enzyme-coding
90  genes and their GSL profiles exhibited across 72 different *A. thaliana* ecotypes. The selection
91  of 72 ecotypes is based on the availability of information about the gene sequences and the
92  patterns of accumulation of aliphatic glucosinolates. Importantly, the experiments for

3

93    determining the aliphatic glucosinolate levels have been performed under identical conditions
94    (Chan, Rowe, & Kliebenstein, 2010; Kliebenstein, Kroymann, et al., 2001). It can therefore
95    be assumed that the environment was identical (up to experimental precision) for all
96    ecotypes. This study presents an effort to quantify the impact of the diversity of MAM gene
97    sequence on GSL variation, rather than the on/off nature of the GS ELONG QTLs. The
98    genetic sequence coding for an enzyme determines the kinetic properties of the corresponding
99    enzyme. For example, polymorphisms in the active sites, in principle, can change the
100   substrate specificity of the respective enzyme. Thus, we investigate the level of
101   polymorphisms in the coding region of the GSL enzymes to study the impact on the diversity
102   of aliphatic GSLs.

## Results

### Distribution of MAM genes across Arabidopsis thaliana ecotypes

105   The genetic basis of chain-length distribution of aliphatic GSLs became evident with the
106   identification of the GS-ELONG QTL in Arabidopsis and *Brassica napus* (Magrath et al.,
107   1994). The locus was mapped in Arabidopsis by using a cross between two ecotypes,
108   *Columbia (Col)* and *Landsberg erecta (Ler)*, where the major glucosinolates are *4C* and *3C*,
109   respectively (Kroymann et al., 2001). The underlying candidates *MAM1* and *MAM3* genes are
110   two adjacent sequences with high similarity to genes encoding isopropylmalate synthase that
111   catalyses the condensation of chain elongation in leucine biosynthesis. Later, a third MAM-
112   like gene, *MAM2*, was identified at the same locus as *MAM1* (Kroymann et al., 2003). While
113   *MAM3* is ubiquitous, most Arabidopsis ecotypes examined possessed functional copies of
114   either *MAM1* or *MAM2* genes. A functional *MAM1* gene has been correlated with the
115   accumulation of 4C GSLs, whereas a functional *MAM2* has been linked to 3C GSLs. To gain
116   more insights on the impact of MAM synthases on chain lengths distribution of aliphatic
117   GSL, we analysed the similarity of the annotated *MAM1* gene across 72 Arabidopsis ecotypes
118   taken from 1001 genome project database (Jorge et al., 2016). These ecotypes were selected
119   based on the availability of gene sequences and the associated aliphatic GSL profiles. Details
120   on the 72 ecotypes are given in the Supplementary Table 1. The annotated *MAM1* sequences
121   from the 72 ecotypes were compared to each other for diversity. Figure 1 shows a mid-point
122   rooted phylogenetic tree showing the evolutionary relationship between the ecotypes based
123   on the similarities and differences in the coding region of the MAM1/MAM2 sequences.
124   Based on maximum likelihood estimation (Guindon et al., 2010), the tree shows two main

125  branches. While 53 out of the 72 ecotypes clustering in the blue branch indeed possess high

126  similarity to the coding region of *MAM1* gene, 19 ecotypes in the red branch possess genes

127  more like the *MAM2* gene. Thus, we assume that the ecotypes composed in blue and red

128  branches possess *MAM1* and *MAM2* genes, respectively.

## The metabolic genotypes and associated phenotypes

129

130  We define a metabolic genotype ($G_i$) as the gene sequence of enzymes of glucosinolate

131  biosynthesis in ecotype $i$, whereas the metabolic phenotype ($P_i$) corresponds to the

132  composition of aliphatic GSLs in the ecotype $i$. To gain a deeper understanding of how

133  different metabolic genotypes and their associated phenotypes are linked, we analysed the

134  genotypic and phenotypic distances. The genotypic distances between the genotypes were

135  calculated as Hamming distance (Hamming, 1950) $d_{i,j}{}^{G}$, which is the number of positions at

136  which the corresponding nucleotide/amino-acid characters are different between gene

137  sequences $G_i$ and $G_j$ of equal length. The phenotypic distance $d_{i,j}{}^{P}$ was calculated as

138  Euclidean distance, $d_{i,j}^{P} = \sqrt{\sum_{k=1}^{n}(P_{i,k} - P_{j,k})^2}$ between two phenotypes $P_{i,k} =$

139  $(P_{i,1}, P_{i,2}, \dots, P_{i,n})$ and $P_{j,k} = (P_{j,1}, P_{j,2}, \dots, P_{j,n})$ in an $n$-dimensional space. In this study,

140  n=6, which corresponds to the total number of chain-elongated aliphatic GSLs found in *A.*

141  *thaliana* (Figure 2). Thus, we can quantify differences between all pairs of ecotypes based on

142  their metabolic genotype $G_i$ ($i = 1, \dots, 72$) and phenotype $P_i$ ($i = 1, \dots, 72$). Figure 3

143  showcases the summary of the analysis, where the genotypic distances are plotted against the

144  phenotypic distances. Intuitively, one would assume that similar genotypes shall exhibit

145  similar phenotypes, and *vice-versa*. However, Figure 3 clearly shows that several ecotypes

146  show low genotypic distance (i.e. they are genotypically similar) but exhibit high phenotypic

147  distance (variation in GSL profiles). Also, there exist genotypically diverse ecotypes that

148  exhibit very similar GSL profiles. Thus, investigating the factors affecting variations in the

149  phenotypes of such ecotypes will provide a clearer understanding of how distinct patterns of

150  GSL accumulation emerge out of genetic differences. Moreover, investigation of the

151  localisation of polymorphic residues in the GSL biosynthesis enzymes will provide a better

152  understanding of the link between metabolic genotype and the associated metabolic

153  phenotypes.

## Diversity of GSL enzyme-coding genes

154

155 To investigate the diversity of metabolic genes of GSL biosynthesis, we investigated the

156 levels of amino acid and nucleotide polymorphism across the 72 *A. thaliana* ecotypes by

157 calculating the average Shannon entropy (Shannon) *H* across the gene length (Figure 4A and

158 B). The analysis revealed that some of these enzymes are highly diverse while others remain

159 conserved across different ecotypes. Interestingly, the diversity seems to be independent of

160 steps of GSL biosynthesis in which the enzymes are active. From the diversity of amino acid

161 sequences (Figure 4A), *FMO-GSOX1* enzyme exhibits the highest diversity (entropy), while

162 the lowest diversity is found in *SOT17*. Among the enzymes active in the chain-elongation

163 pathway, *MAM1* shows the highest diversity while *BAT5* shows the low diversity. However, a

164 further analysis of the diversity in the nucleotide sequences of the metabolic genotypes

165 showed a high level of polymorphism in BAT5 (cf. Figure 4B), which was not reflected in

166 the diversity of amino acid sequences. Indeed, most genes show only a slightly lower

167 diversity in amino acid variation than nucleotide variation, which reflects the degeneration of

168 the genetic code (Figure 4C). A plausible explanation for the low amino acid variation in

169 BAT5 could be the specificity towards a variety of chain-elongated substrates of GSL

170 biosynthesis (Halkier & Gershenzon, 2006). The low diversity of BAT5 could be linked to its

171 function in transport of a diverse range of compounds that are a part of aliphatic GSL

172 biosynthesis and Met-salvage pathway (Gigolashvili et al., 2009; Sauter, Moffatt, Saechao,

173 Hell, & Wirtz, 2013), thus, mutations in the coding region of BAT5 may impair the

174 functioning of both pathways. In contrast to *BAT5*, *FMO-GSOX1* shows a low diversity in the

175 nucleotide sequences of 72 genotypes but reflect a high diversity in the amino acid

176 sequences. This is a clear example of preferential accumulation of non-synonymous

177 mutations, which alter the amino acid sequence of an enzyme.

178 High diversity of *MAM1* could be a consequence of incorrect annotation of *MAM1*/*MAM2*

179 enzymes across 72 *A. thaliana* ecotypes. Thus, we analysed the diversity of GSL enzymes

180 across the MAM1 ecotypes (blue branch of Figure 1) and MAM2 (red branch of Figure 1)

181 ecotypes, separately. We did see a reduction in the diversity of MAM1 and MAM2 enzymes

182 (cf. Supplementary Figure 1 and Supplementary Figure 2). Nevertheless, *MAM1* and *MAM2*

183 are still the most diverse enzymes of chain-elongation pathway.

## Polymorphisms in the active sites of MAM enzymes

185 To get a clearer understanding of the effects of the localisation of polymorphic amino acid

186 residues in the active sites of the metabolic enzymes, we extracted the information about the

187 active sites from the NCBI's conserved domain database (Marchler-Bauer et al., 2015). For

188 example, the amino acid positions 93, 94, 97, 124, 162, 164, 186, 227, 229, 231, 257, 259,

189 260, 261, 262, 290, 292, and 294 are known to be key for activity of MAM synthases, based

190 on the database and a crystal structure (Kumar et al., 2019; Marchler-Bauer et al., 2015;

191 Petersen et al., 2019). We refer to the amino acid positions from 93 to 294 as active region of

192 the enzyme. We have found that MAM synthases exhibit a maximum of 13 and 3

193 polymorphic residues in the active region of MAM1 and MAM3, respectively. Figure 5(A)

194 and (B) show pairwise comparisons of polymorphisms in the active region of MAM

195 synthases against the genotypic distances of 72 *A. thaliana* ecotypes. Furthermore, we

196 recorded the polymorphisms at different positions in the active region of MAM synthases

197 (see Figure 5(C)). The amino acid residues at position 98, 99, 132, 138, 139, 147, 165, 173,

198 177, 187, 228, 245, 257, 258, 271, 289 and 290 of the *MAM1* and positions 156, 231, 241 and

199 242 of *MAM3* accumulate polymorphic residues across the 72 Arabidopsis ecotypes.

200 Polymorphisms in the active sites of an enzyme, in principle, can change the catalytic

201 properties of the enzyme (Kroymann et al., 2001; Kumar et al., 2019; Petersen et al., 2019).

202 However, the quantitative effect on the enzymatic properties cannot be explained due to

203 unavailability of enzyme abundances in these 72 ecotypes.

## Discussion

### Plasticity of the metabolic genotype and the associated GSL profiles

206 Glucosinolate metabolism results in a highly variable composition of individual metabolites

207 in Arabidopsis accessions, which is reflected by a corresponding high diversity at the causal

208 genetic level. Thus, it serves as a suitable model system to investigate the broader aspects of

209 genotype-phenotype relationships. Allelic composition at several glucosinolate biosynthetic

210 loci drive different glucosinolate profiles among *A. thaliana* ecotypes (Kliebenstein,

211 Kroymann, et al., 2001). These variations, however, are often in the form of presence and/or

212 expression of one or other copy of a duplicated gene, such as the *AOP2* and *AOP3* at the GS-

213 ALK/GS-OHP locus (Kliebenstein, Lambrix, Reichelt, Gershenzon, & Mitchell-Olds, 2001),

214 or the *MAM1*/*MAM2* at GS-ELONG (Kroymann et al., 2003), which complicates our

215 understanding of how genetic variations lead to metabolic properties of the enzymes encoded

216 by the respective genes. The Met-derived aliphatic GSLs are the most abundant form of

217 glucosinolates in *A. thaliana* and many *Brassicaceae* crops (Agerbirk & Olsen, 2012;

218 Benderoth et al., 2006; Halkier & Gershenzon, 2006; Kliebenstein, Kroymann, et al., 2001;

219     Kroymann et al., 2003; Textor et al., 2004; Textor, de Kraker, Hause, Gershenzon, &

220     Tokuhisa, 2007). The chain-elongation pathway of GSL biosynthesis is crucial for generating

221     the chain-length diversity of aliphatic GSLs and for connecting primary and specialised

222     metabolism. Although the evolution of core features of aliphatic GSL biosynthesis in

223     Arabidopsis has been studied (He et al., 2011; Sawada et al., 2009; Textor et al., 2004;

224     Wittstock et al., 2004), the molecular basis for diversity of function of MAM synthases and

225     the role of different MAM isoforms within *A. thaliana* accessions is not sufficiently

226     understood. The chain-length distribution of different aliphatic GSLs has so far been

227     attributed to the presence of different MAM isoforms, namely *MAM1*, *MAM2* and *MAM3*

228     (Halkier & Gershenzon, 2006; Kroymann et al., 2003; Textor et al., 2004, 2007). By

229     investigating the differences in sequences of MAM synthases across *A. thaliana* ecotypes, we

230     show that the ecotypes can be broadly classified in two groups based on the similarity to

231     either of *MAM1* and *MAM2* genes as also expected from the genomic composition of the

232     *GS_ELONG* locus (cf. Figure 1). Correspondingly, the GSL profiles from different ecotypes

233     can be broadly classified into two major groups based on the phenotypic distance $d_{i,j}^{P}$

234     between different GSL profiles (cf. Figure 2). However, the groups classified based on either

235     of phenotypic distances or the similarity of MAM synthases are not identical. This points to a

236     more complicated relationship between the genotype and the associated GSL profiles. Thus,

237     estimating the pattern of GSL accumulation based solely on the distinction between MAM1

238     and MAM2 enzymes is not feasible.

239     ## Diversity of genes beyond classical QTLs

240     The heterogeneity in the genetic makeup of the metabolic genes across different *A. thaliana*

241     ecotypes and their associated metabolic phenotypes are an excellent tool for investigating the

242     mechanisms of adaptation and functional diversification (Mitchell-Olds & Schmitt, 2006;

243     Pigliucci, 2010). Comparing the diversity of genes of GSLs synthesis we expected to find the

244     highest diversity in genes of the chain-elongation and secondary modification, because these

245     steps contribute highly to the diversity of the GSLs. However, surprisingly, we found that the

246     level of diversity appears to be unrelated with the functional role of the gene within the GSL

247     metabolic pathway (cf. Figure 4). While as expected, the least diverse enzyme was SOT17,

248     part of the core synthesis, another enzyme of the core pathway, SUR1, was the fourth most

249     diverse from 30 enzymes (Figure 4). This is surprising, since loss of enzymes of the core

250     synthesis, such as *UGT74B1* or *SUR1* has a much higher impact on the total GSLs than loss

251     of enzymes of the side chain modification (Douglas Grubb et al., 2004; Keurentjes et al.,

252    2006; Mikkelsen, Naur, & Halkier, 2004). Thus, it seems that enzymes of all parts of GSL

253    biosynthesis can contribute to the diversity of the metabolites. Surprisingly, while

254    investigating the second least diverse enzyme *BAT5*, part of chain-elongation of aliphatic

255    GSLs, we found that it is highly diverse in the nucleotide sequence (see Figure 4B). It can be

256    a consequence of purifying natural selection that prevents the change of an amino acid

257    residue at a given position in a multiple alignment, thus favouring an excess of synonymous

258    versus non-synonymous substitutions. It is much more difficult, however, to explain the

259    results of *FMO-GSOX1*, which shows a large variation in amino acid sequence derived from

260    a relatively low variation in DNA sequence.  Nevertheless, it is evident that *FMO-GSOX1*

261    favours non-synonymous substitutions versus the synonymous substitutions, possibly linked

262    to the function in the secondary modification of glucosinolates, responsible for large part of

263    the structural variation. Multiple genetic analyses revealed that, in general, few key QTLs

264    shape the metabolic phenotype. In contrast, our analysis detected diversity of the metabolic

265    genes across the whole pathway, irrespective of their association to a major QTL. However,

266    how far this variation in gene/protein sequence contributes to the phenotypic variation

267    remains to be elucidated.

268    ## The genotype-phenotype relationship

269    Nowhere is the contribution of subtle sequence diversity to variation in GSLs more apparent,

270    than in the MAM genes. In *A. thaliana*, the enzyme isoforms *MAM1* and *MAM2* catalyse the

271    formation of short-chain (*3C* and *4C*) aliphatic GSLs, whereas the isoform MAM3 catalyses

272    the formation of both short-chain and long-chain (*5C-8C*) GSLs (Halkier2006). Indeed,

273    orthologues of *MAM1* and *MAM2* are also responsible for diversity of aliphatic GSLs across a

274    range of *Brassica* species (Kumar et al., 2019). The distinct function of the two genes was

275    confirmed by complementation of Arabidopsis *mam1* mutant, when *MAM1* from *B. juncea*

276    restored wild type GSL profile but *MAM2* did not (Kumar et al., 2019). Our investigation of

277    genotypic and phenotypic distances between different Arabidopsis ecotypes showed that

278    some ecotypes have identical metabolic genotype but exhibit high diversity in their associated

279    metabolic phenotype, and *vice-versa* (cf. Figure 3). Therefore, the relationship between the

280    metabolic genotypes and the phenotypes are much more complicated than a link to one of the

281    *MAM1*/*MAM2* gene pair. Indeed, the role of individual amino acid alterations between these

282    two genes demonstrates clearly that also SNPs can have a great effect on the phenotypes.

283    Thus, mutagenesis of serine to phenylalanine at position 102 and alanine to threonine at 290,

284    parts of active region of MAM1 changed the distribution of C3 to C4 GSLs (Kroymann et al.,

285  2001) in *A. thaliana*. Alterations in four other amino acids in *B. juncea* MAM1 affected the

286  kinetic properties of the enzyme to more MAM2-like and *vice versa* (Kumar et al., 2019).

287  Also, in Arabidopsis MAM1 alteration of further three amino acids resulted in changes of the

288  pattern of elongation products *in vitro* (Petersen et al., 2019). To further investigate the

289  impact of MAM synthases on chain-length distribution of aliphatic GSLs, we analysed the

290  polymorphisms in the active region of MAM synthases. From our analysis of the diversity of

291  the active region of MAM synthases, we conclude that *MAM1* is highly variable across its

292  active region and accumulate up to 13 polymorphic amino-acid residues at 17 different

293  locations in the active region. Whereas, the active region of *MAM3* is comparatively

294  conserved and only accumulates a maximum of 4 polymorphic residues at 3 locations in the

295  active region, naturally (cf. Figure 5). It is known that polymorphisms in active site of MAM

296  synthases change the specificity of a metabolic enzyme towards respective substrates of

297  aliphatic GSL biosynthesis (Kumar et al., 2019; Petersen et al., 2019). This results in

298  different composition and the total accumulation of aliphatic GSLs across different

299  *Brassicaceae* species, including *Arabidopsis thaliana* (Kroymann et al., 2001; Kumar et al.,

300  2019; Petersen et al., 2019). The above analysis, however, only showcases one of the two

301  possibilities by which a genotype can exhibit a metabolic phenotype. The gene regulatory

302  networks can also change the expression of metabolic genes, which in turn changes the

303  enzyme abundance and thus results in different metabolic phenotypes (de Kraker &

304  Gershenzon, 2011; Kumar et al., 2019; Petersen et al., 2019). Although, numerous studies

305  have shown that a multitude of genes and underlying regulatory processes are involved in the

306  diversity of specialised metabolites such as glucosinolates (Chan et al., 2010; Koornneef,

307  Alonso-Blanco, & Vreugdenhil, 2004; Kumar et al., 2019; LASKY et al., 2012; Petersen et

308  al., 2019), interpreting the findings in the context of metabolic properties is highly

309  challenging. This is particularly due to a missing stringent definition of the genotype–

310  phenotype relationship, which can hardly be expected to be derivable from a single

311  methodology but rather requires a comprehensive platform of combined experimental and

312  theoretical strategies (DIZ, MARTÍNEZ-FERNÁNDEZ, & ROLÁN-ALVAREZ, 2012;

313  Sharma, 2018; Weckwerth, Wenzel, & Fiehn, 2004).

## Conclusion
314

315  Altogether we show here that the control over phenotypic diversity in glucosinolates is

316  potentially spread over the whole pathway. On the example of MAM1 and MAM2,

317  responsible for side chain elongation of Met-derived glucosinolates, we revealed that

318    sequence variation beyond the presence of one or the other isoform contributes to the

319    variation in chain length. The present study thus points to the necessity to pay attention to

320    variation beyond the classical ON/OFF features of key metabolic QTLs, for investigating the

321    diversity of specialised metabolic pathways, such as glucosinolates. Since the recent efforts

322    towards unravelling the genotype-phenotype relationships focus at either experimental

323    studies with a selection of genotypes or computational approaches to correlate the observed

324    experimental observations, it is crucial to develop frameworks that integrate multi-omics data

325    with fundamental rules of metabolic modelling to fully understand how particular genotype is

326    reflected in a phenotype.

## Materials and Methods

327

### Genotypic data

328

329    Information about the nucleotide and amino acid composition of 30 GSL biosynthesis genes

330    from 72 *A. thaliana* ecotypes was taken from the 1001 genomes project (Jorge et al., 2016).

331    The reason behind selecting these 72 ecotypes was the availability of Met-derived GSL

332    composition under identical environmental conditions. To obtain the gene sequences of the

333    ecotypes of interest, we used an inhouse R-script that converts the TAIR10 version of SNP

334    (single nucleotide polymorphisms) files provided by the 1001 genome database into an R-

335    object. This R-object is a sparse matrix containing the nucleotide information for each

336    ecotype at each locus in the reference genome coded as numbers. Non-polymorphic sites are

337    coded as 0, polymorphic ones as 1,2,3,4 or 5 depending if at the specific locus an A, C, G, T

338    or indel was observed. From this R-Object we could extract the nucleotide sequences of a

339    specific ecotype for each coding region of interest. To obtain the amino acid sequence we

340    used the function *'translate'* from the R-package 'seqinr'.

### Phenotypic data

341

342    Experimental data of Met-GSLs concentrations were obtained from Chan et al. (2010) and

343    Kliebenstein et al. (2001). The final set of GSL data is given in Table 1. The data is

344    composed of normalised concentrations of six aliphatic GSLs, referred to as 3C to 8C, from

345    72 different ecotypes of *A. thaliana* under controlled and identical experimental conditions.

346

### Calculation of diversity using Shannon entropy

347

348    The nucleotide/amino-acid composition of the coding regions of GSL genes from 72 *A.*

349    *thaliana* ecotypes is described as a set of relative probability, $p_{i,j}$, for the $i^{th}$ nucleotide/amino-

350    acid ($i = 1, 2, ..., n$) in the $j^{th}$ ecotype ($j = 1, 2, ..., 72$). Then, the diversity of each position

351    in the coding region can be quantified by Shannon entropy (Shannon, 1948),

$$H_i = -\sum_{j}^{n} p_{i,j} \log p_{i,j}$$

352    $H_i$ will vary from zero, when the $i^{th}$ nucleotide/amino-acid is same across all 72 ecotypes, to 1

353    when the probability is equal for observing all nucleotides/amino-acids at same locus across

354    72 ecotypes. Moreover, to get an estimate of diversity of a nucleotide/amino-acid sequence of

355    length *n*, we calculate the average entropy $H^{avg}$ as

$$H^{avg} = \frac{1}{n} \sum_{i}^{n} H_i$$

## Phylogenetic tree reconstruction

357    Amino acid sequences of the MAM loci from the 72 ecotypes were aligned using 'mafft' ver.

358    v7.407 (Katoh & Standley, 2013) with the parameter '--maxiterate 1000 --globalpair --

359    phylipout' to obtain a multiple sequences alignment in phylip format. This was then used as

360    input for phyml (20120412) (Guindon et al., 2010) to reconstruct the maximum likelihood

361    phylogenetic tree using the LG substitution model. The tree was visualised with FigTree

362    v1.4.2 (http://tree.bio.ed.ac.uk/software/figtree/).

## Resources

364    The data and Python scripts used to produce the results presented in this manuscript are

365    available with instructions at (https://gitlab.com/surajsept/GTvsPT).

## Acknowledgements

## Author contributions

371    SS, OE and SK planned and designed the research. SS performed the computational work and

372    wrote the manuscript with the help of OP, OE and SK.

# References

Agerbirk, N., & Olsen, C. E. (2012). Glucosinolate structures in evolution. *Phytochemistry*, *77*, 16–45. https://doi.org/10.1016/j.phytochem.2012.02.005

Agerbirk, N., Olsen, C. E., Heimes, C., Christensen, S., Bak, S., & Hauser, T. P. (2015). Multiple hydroxyphenethyl glucosinolate isomers and their tandem mass spectrometric distinction in a geographically structured polymorphism in the crucifer Barbarea vulgaris. *Phytochemistry*, *115*, 130–142. https://doi.org/10.1016/J.PHYTOCHEM.2014.09.003

Benderoth, M., Textor, S., Windsor, A. J., Mitchell-Olds, T., Gershenzon, J., & Kroymann, J. (2006). Positive selection driving diversification in plant secondary metabolism. *Proceedings of the National Academy of Sciences of the United States of America*, *103*(24), 9118–9123. https://doi.org/10.1073/pnas.0601738103

Buskov, S., Serra, B., Rosa, E., Sørensen, H., & Sørensen, J. C. (2002). Effects of intact glucosinolates and products produced from glucosinolates in myrosinase-catalyzed hydrolysis on the potato cyst nematode (Globodera rostochiensis Cv. Woll). *Journal of Agricultural and Food Chemistry*, *50*(4), 690–695. Retrieved from http://www.ncbi.nlm.nih.gov/pubmed/11829629

Chan, E. K. F., Rowe, H. C., & Kliebenstein, D. J. (2010). Understanding the evolution of defense metabolites in Arabidopsis thaliana using genome-wide association mapping. *Genetics*, *185*(3), 991–1007. https://doi.org/10.1534/genetics.109.108522

Clarke, D. B. (2010). Glucosinolates, structures and analysis in food. *Analytical Methods*, *2*(4), 310. https://doi.org/10.1039/b9ay00280d

Daxenbichler, M. E., Spencer, G. F., Carlson, D. G., Rose, G. B., Brinker, A. M., & Powell, R. G. (1991). Glucosinolate composition of seeds from 297 species of wild plants. *Phytochemistry*, *30*(8), 2623–2638.

de Kraker, J.-W., & Gershenzon, J. (2011). From amino acid to glucosinolate biosynthesis: protein sequence changes in the evolution of methylthioalkylmalate synthase in Arabidopsis. *The Plant Cell*, *23*(1), 38–53. https://doi.org/10.1105/tpc.110.079269

DIZ, A. P., MARTÍNEZ-FERNÁNDEZ, M., & ROLÁN-ALVAREZ, E. (2012). Proteomics in evolutionary ecology: linking the genotype with the phenotype. *Molecular Ecology*,

403    *21*(5), 1060–1080. https://doi.org/10.1111/j.1365-294X.2011.05426.x

404    Douglas Grubb, C., Zipp, B. J., Ludwig-Müller, J., Masuno, M. N., Molinski, T. F., & Abel,
405        S. (2004). Arabidopsis glucosyltransferase UGT74B1 functions in glucosinolate
406        biosynthesis and auxin homeostasis. *The Plant Journal*, *40*(6), 893–908.

407    Fahey, J. W., Zalcmann,  a T., & Talalay, P. (2001). The chemical diversity and distribution
408        of glucosinolates and isothiocyanates amoung plants. *Phytochemistry*, *56*, 5–51.

409    Fisher, R. A. (1930). GENETICAL THEORY OF NATURAL SELECTION. Retrieved from
410        http://14.139.56.90/bitstream/1/2033620/1/IVRI 3205.pdf

411    Gabrys, B., & Tjallingii, W. F. (2002). The role of sinigrin in host plant recognition by aphids
412        during initial plant penetration. *Entomologia Experimentalis et Applicata*, *104*(1), 89–
413        93. https://doi.org/10.1046/j.1570-7458.2002.00994.x

414    Gigolashvili, T., Yatusevich, R., Rollwitz, I., Humphry, M., Gershenzon, J., & Flügge, U.-I.
415        (2009). The plastidic bile acid transporter 5 is required for the biosynthesis of
416        methionine-derived glucosinolates in Arabidopsis thaliana. *The Plant Cell*, *21*(6), 1813–
417        1829. https://doi.org/10.1105/tpc.109.066399

418    Guindon, S., Dufayard, J.-F., Lefort, V., Anisimova, M., Hordijk, W., & Gascuel, O. (2010).
419        New algorithms and methods to estimate maximum-likelihood phylogenies: assessing
420        the performance of PhyML 3.0. *Systematic Biology*, *59*(3), 307–321.

421    Halkier, B. A., & Gershenzon, J. (2006). BIOLOGY AND BIOCHEMISTRY OF
422        GLUCOSINOLATES. *Annual Review of Plant Biology*, *57*(1), 303–333.
423        https://doi.org/10.1146/annurev.arplant.57.032905.105228

424    Hamming, R. W. (1950). Error detecting and error correcting codes. *Bell Labs Technical*
425        *Journal*, *29*(2), 147–160.

426    He, Y., Galant, A., Pang, Q., Strul, J. M., Balogun, S. F., Jez, J. M., & Chen, S. (2011).
427        Structural and functional evolution of isopropylmalate dehydrogenases in the leucine
428        and glucosinolate pathways of Arabidopsis thaliana. *Journal of Biological Chemistry*,
429        *286*(33), 28794–28801. https://doi.org/10.1074/jbc.M111.262519

430    Jorge, C., Becker, C., Bemm, F., Bergelson, J., Borgwardt, K. M., Cao, J., … Zhou, X.
431        (2016). 1,135 Genomes Reveal the Global Pattern of Polymorphism in Arabidopsis
432        thaliana. *Cell*, *166*(2), 481–491. https://doi.org/10.1016/j.cell.2016.05.063

433    Katoh, K., & Standley, D. M. (2013). MAFFT Multiple Sequence Alignment Software
434        Version 7: Improvements in Performance and Usability. *Molecular Biology and*
435        *Evolution*, *30*(4), 772–780. https://doi.org/10.1093/molbev/mst010

436    Keurentjes, J. J. B., Fu, J., de Vos, C. H. R., Lommen, A., Hall, R. D., Bino, R. J., …
437        Koornneef, M. (2006). The genetics of plant metabolism. *Nature Genetics*, *38*(7), 842–
438        849. https://doi.org/10.1038/ng1815

439    Kliebenstein, D. J., Kroymann, J., Brown, P., Figuth, A., Pedersen, D., Gershenzon, J., &
440        Mitchell-Olds, T. (2001). Genetic control of natural variation in Arabidopsis
441        glucosinolate accumulation. *Plant Physiology*, *126*(2), 811–825.
442        https://doi.org/10.1104/pp.126.2.811

443    Kliebenstein, D. J., Lambrix, V. M., Reichelt, M., Gershenzon, J., & Mitchell-Olds, T.
444        (2001). Gene Duplication in the Diversification of Secondary Metabolism: Tandem 2-
445        Oxoglutarate-Dependent Dioxygenases Control Glucosinolate Biosynthesis in
446        Arabidopsis. *THE PLANT CELL ONLINE*, *13*(3), 681–693.
447        https://doi.org/10.1105/tpc.13.3.681

448    Koornneef, M., Alonso-Blanco, C., & Vreugdenhil, D. (2004). NATURALLY OCCURRING
449        GENETIC VARIATION IN *ARABIDOPSIS THALIANA*. *Annual Review of Plant*
450        *Biology*, *55*(1), 141–172. https://doi.org/10.1146/annurev.arplant.55.031903.141605

451    Kroymann, J., Donnerhacke, S., Schnabelrauch, D., & Mitchell-Olds, T. (2003). Evolutionary
452        dynamics of an Arabidopsis insect resistance quantitative trait locus. *Proceedings of the*
453        *National Academy of Sciences of the United States of America*, *100 Suppl*(Supplement
454        2), 14587–14592. https://doi.org/10.1073/pnas.1734046100

455    Kroymann, J., Textor, S., Tokuhisa, J. G., Falk, K. L., Bartram, S., Gershenzon, J., &
456        Mitchell-Olds, T. (2001). A gene controlling variation in Arabidopsis glucosinolate
457        composition is part of the methionine chain elongation pathway. *Plant Physiology*,
458        *127*(3), 1077–1088. https://doi.org/10.1104/pp.010416

459    Kumar, R., Lee, S. G., Augustine, R., Reichelt, M., Vassão, D. G., Palavalli, M. H., … Bisht,
460        N. C. (2019). Molecular Basis of the Evolution of Methylthioalkylmalate Synthase and
461        the Diversity of Methionine-Derived Glucosinolates. *The Plant Cell*, *31*(7), 1633–1647.
462        https://doi.org/10.1105/tpc.19.00046

463    LASKY, J. R., DES MARAIS, D. L., McKAY, J. K., RICHARDS, J. H., JUENGER, T. E.,

464       & KEITT, T. H. (2012). Characterizing genomic variation of *Arabidopsis thaliana*□: the

465       roles of geography and climate. *Molecular Ecology*, *21*(22), 5512–5529.

466       https://doi.org/10.1111/j.1365-294X.2012.05709.x

467    Lazzeri, L., Curto, G., Leoni, O., & Dallavalle, E. (2004). Effects of Glucosinolates and Their

468       Enzymatic Hydrolysis Products via Myrosinase on the Root-knot Nematode

469       *Meloidogyne incognita* (Kofoid et White) Chitw. *Journal of Agricultural and Food*

470       *Chemistry*, *52*(22), 6703–6707. https://doi.org/10.1021/jf030776u

471    Lynch, M., Walsh, B., & others. (1998). *Genetics and analysis of quantitative traits* (Vol. 1).

472       Sinauer Sunderland, MA.

473    Magrath, R., Banot, F., Morgner, M., Parkin, I., Sharpe, A., Lister, C., … Mithen, A. A.

474       (1994). Genetical Society of Great Britain Genetics of aliphatic glucosinolates. I. Side

475       chain elongation in Brassica napus and A rabidopsis thaliana. *Heredity*, *72*, 290–299.

476       Retrieved from https://www.jic.ac.uk/staff/caroline-

477       dean/pdf_files/88_Magrath_R_et_al_1994_Heredity.pdf

478    Marchler-Bauer, A., Derbyshire, M. K., Gonzales, N. R., Lu, S., Chitsaz, F., Geer, L. Y., …

479       Bryant, S. H. (2015). CDD: NCBI's conserved domain database. *Nucleic Acids*

480       *Research*, *43*(Database issue), D222. https://doi.org/10.1093/NAR/GKU1221

481    Mewis, I., Ulrich, C., & Schnitzler, W. H. (2002). The role of glucosinolates and their

482       hydrolysis products in oviposition and host-plant finding by cabbage webworm, Hellula

483       undalis. *Entomologia Experimentalis et Applicata*, *105*(2), 129–139.

484       https://doi.org/10.1046/j.1570-7458.2002.01041.x

485    Mikkelsen, M. D., Naur, P., & Halkier, B. A. (2004). Arabidopsis mutants in the C--S lyase

486       of glucosinolate biosynthesis establish a critical role for indole-3-acetaldoxime in auxin

487       homeostasis. *The Plant Journal*, *37*(5), 770–777.

488    Miles, C. I., Campo, M. L. del, & Renwick, J. A. A. (2005). Behavioral and chemosensory

489       responses to a host recognition cue by larvae of Pieris rapae. *Journal of Comparative*

490       *Physiology A*, *191*(2), 147–155. https://doi.org/10.1007/s00359-004-0580-x

491    Mitchell-Olds, T., & Schmitt, J. (2006). Genetic mechanisms and evolutionary significance

492       of natural variation in Arabidopsis. *Nature*, *441*(7096), 947–952.

493    https://doi.org/10.1038/nature04878

494    Noret, N., Meerts, P., Tolrà, R., Poschenrieder, C., Barceló, J., & Escarre, J. (2005).

495    Palatability of Thlaspi caerulescens for snails: influence of zinc and glucosinolates. *New*

496    *Phytologist*, *165*(3), 763–772. https://doi.org/10.1111/j.1469-8137.2004.01286.x

497    Owen, C., Patron, N. J., Huang, A., & Osbourn, A. (2017). Harnessing plant metabolic

498    diversity. *Current Opinion in Chemical Biology*, *40*, 24–30.

499    https://doi.org/10.1016/j.cbpa.2017.04.015

500    Petersen, A., Hansen, L. G., Mirza, N., Crocoll, C., Mirza, O., & Halkier, B. A. (2019).

501    Changing substrate specificity and iteration of amino acid chain elongation in

502    glucosinolate biosynthesis through targeted mutagenesis of Arabidopsis

503    methylthioalkylmalate synthase 1. *Bioscience Reports*, *39*(7).

504    https://doi.org/10.1042/BSR20190446

505    Pigliucci, M. (2010). Genotype–phenotype mapping and the end of the 'genes as blueprint'

506    metaphor. *Philosophical Transactions of the Royal Society B: Biological Sciences*,

507    *365*(1540), 557–566. https://doi.org/10.1098/rstb.2009.0241

508    Rodman, J. E. (1980). Population variation and hybridization in sea-rockets (Cakile,

509    Cruciferae): seed glucosinolate characters. *American Journal of Botany*, 1145–1159.

510    Sauter, M., Moffatt, B., Saechao, M. C. C., Hell, R., & Wirtz, M. (2013). Methionine salvage

511    and S-adenosylmethionine: essential links between sulfur, ethylene and polyamine

512    biosynthesis. *Biochemical Journal*, *451*(2), 145–154.

513    https://doi.org/10.1042/BJ20121744

514    Sawada, Y., Kuwahara, A., Nagano, M., Narisawa, T., Sakata, A., Saito, K., & Yokota Hirai,

515    M. (2009). Omics-based approaches to methionine side chain elongation in Arabidopsis:

516    characterization of the genes encoding methylthioalkylmalate isomerase and

517    methylthioalkylmalate dehydrogenase. *Plant and Cell Physiology*, *50*(7), 1181–1190.

518    Shannon, C. E. (1948). A Mathematical Theory of Communication. *Bell System Technical*

519    *Journal*, *27*(3), 379–423. https://doi.org/10.1002/j.1538-7305.1948.tb01338.x

520    Sharma, S. (2018). Mathematical models of glucosinolate metabolism in plants. Retrieved

521    from https://docserv.uni-duesseldorf.de/servlets/DocumentServlet?id=46153

522    Sønderby, I. E., Geu-Flores, F., & Halkier, B. a. (2010). Biosynthesis of glucosinolates - gene

523    discovery and beyond. *Trends in Plant Science*, *15*(5), 283–290.

524    https://doi.org/10.1016/j.tplants.2010.02.005

525    Textor, S., Bartram, S., Kroymann, J. J. J., Falk, K. L., Hick, A., Pickett, J. a., & Gershenzon,

526    J. (2004). Biosynthesis of methionine-derived glucosinolates in Arabidopsis thaliana ：

527    recombinant expression and characterization of methylthioalkylmalate synthase, the

528    condensing enzyme of the chain-elongation cycle. *Planta*, *218*(6), 1026–1035.

529    https://doi.org/10.1007/s00425-003-1184-3

530    Textor, S., de Kraker, J.-W., Hause, B., Gershenzon, J., & Tokuhisa, J. G. (2007). MAM3

531    catalyzes the formation of all aliphatic glucosinolate chain lengths in Arabidopsis. *Plant*

532    *Physiology*, *144*(1), 60–71. https://doi.org/10.1104/pp.106.091579

533    Traka, M., & Mithen, R. (2009). Glucosinolates, isothiocyanates and human health.

534    *Phytochemistry Reviews*, *8*(1), 269–282. https://doi.org/10.1007/s11101-008-9103-7

535    Weckwerth, W., Wenzel, K., & Fiehn, O. (2004). Process for the integrated extraction,

536    identification and quantification of metabolites, proteins and RNA to reveal their co-

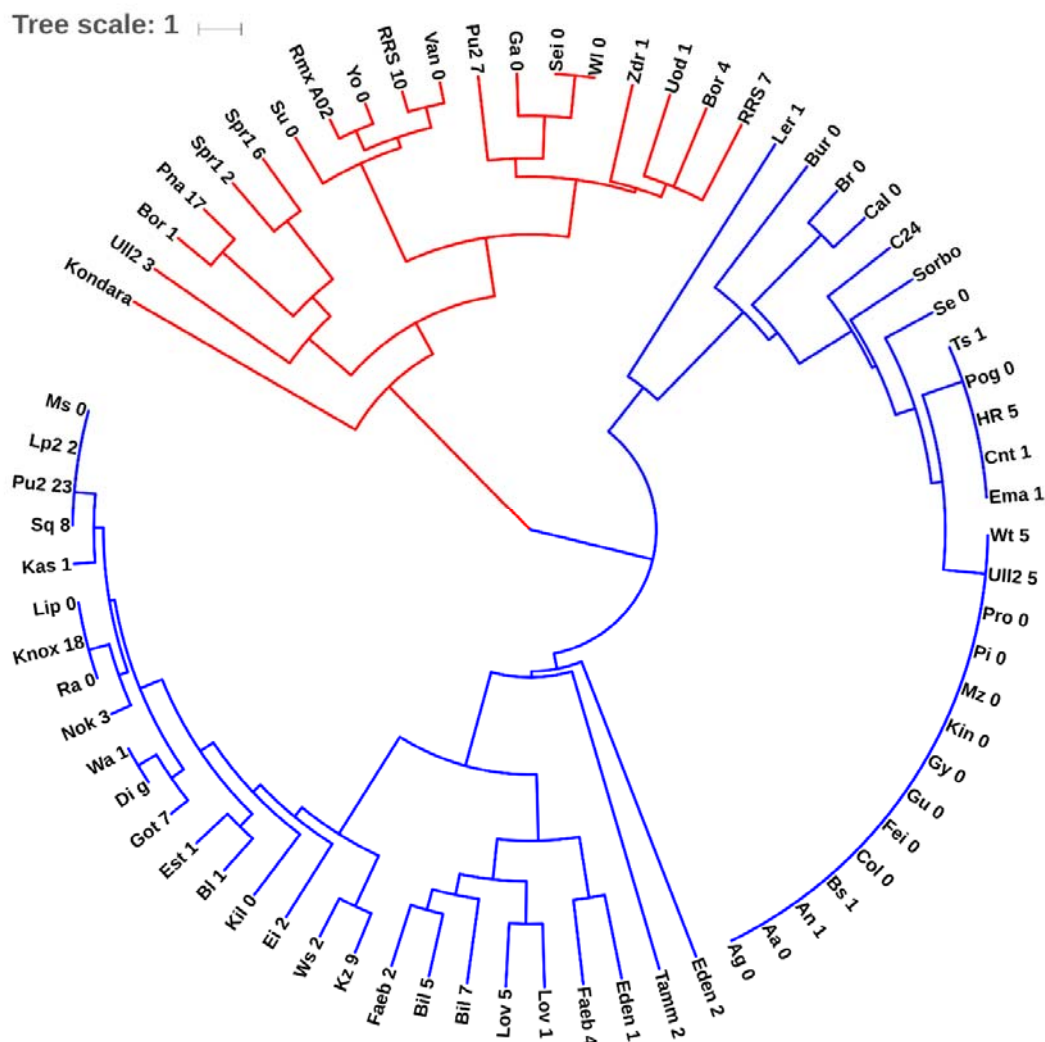537    regulation in biochemical networks. *PROTEOMICS*, *4*(1), 78–83.

538    https://doi.org/10.1002/pmic.200200500

539    Wittstock, U., Agerbirk, N., Stauber, E. J., Olsen, C. E., Hippler, M., Mitchell-Olds, T., …

540    Vogel, H. (2004). Successful herbivore attack due to metabolic diversion of a plant

541    chemical defense. *Proceedings of the National Academy of Sciences*, *101*(14), 4859–

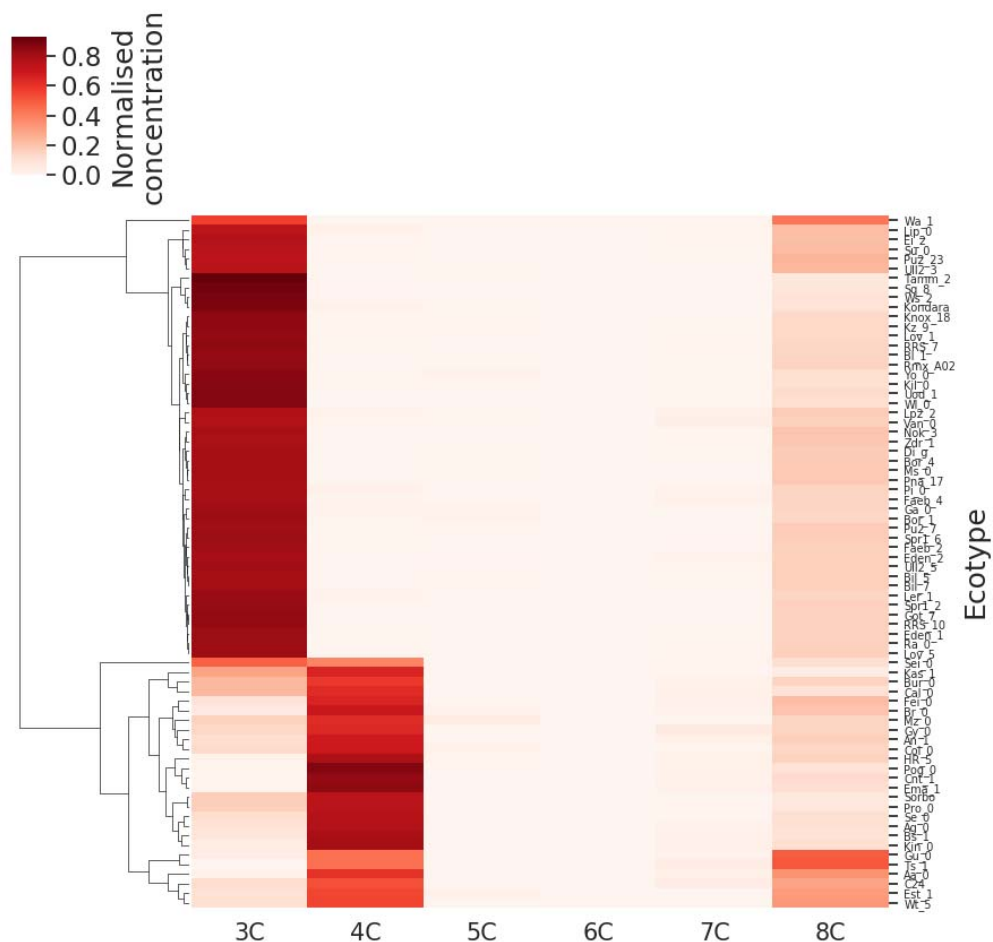542    4864. https://doi.org/10.1073/pnas.0308007101

543

544

# Figures

*Figure 1: Mid-point rooted phylogenetic tree showing the evolutionary relationship among coding regions of annotated MAM1 gene of 72 Arabidopsis thaliana ecotypes. The red branch represents ecotypes showing higher similarity to the MAM2 sequence. The scale bar is substitutions per position.*

552 *Figure 2: Clustered heatmap of the GSL composition across 72 different A. thaliana ecotypes. The top cluster is composed of*
553 *ecotypes having high concentration of 3C GSLs, the lower cluster corresponds to the ecotypes accumulating high*
554 *concentrations of 4C GSLs.*



555

556

557   *Figure 3: The genotypic versus phenotypic distance. Every ecotype is assigned to possess either MAM1 or MAM2, based on*
558   *the sequence similarity of annotated MAM1 gene. Each dot represents a pair of ecotypes. The colours red and blue denote*
559   *pairs, where both ecotypes show high similarity to the coding region of MAM1 and MAM2 sequence, respectively, the green*
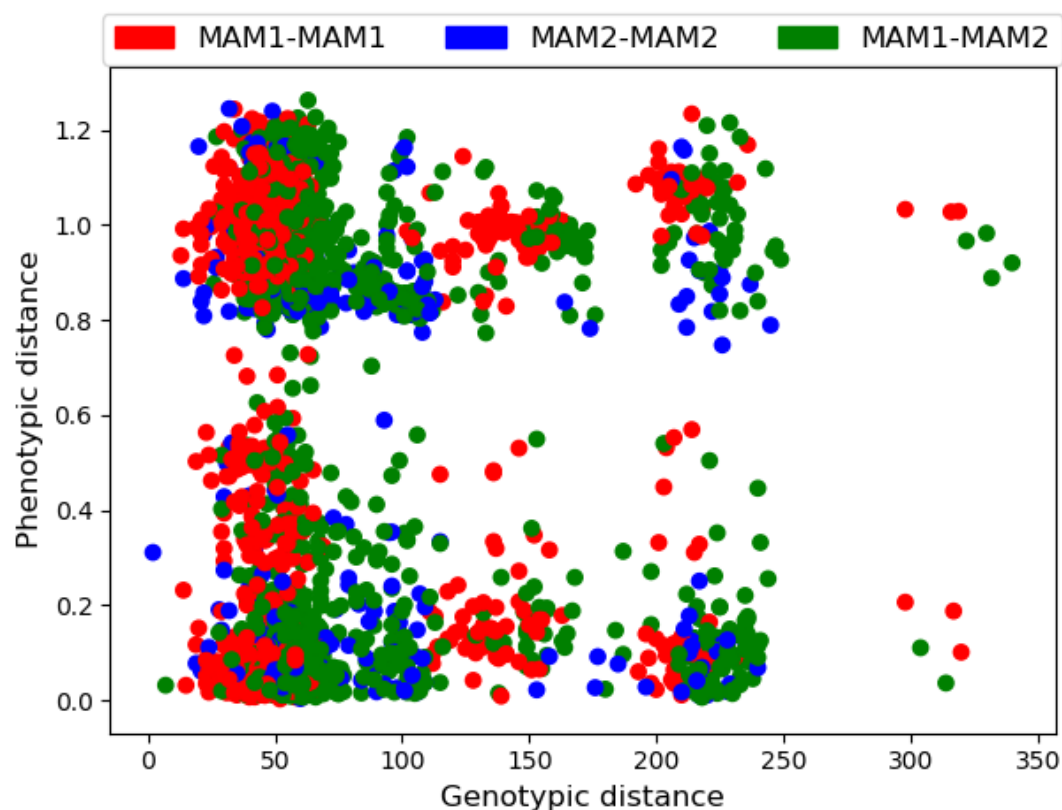560   *dots denote heterogeneous pairs.*



561
562

563 *Figure 4: Diversity of GSL genes from 72 A. thaliana ecotypes. (A) The bars represent the diversity of the amino-acid (AA)*
564 *residues of respective GSL genes. The bars are colour coded to denote genes from the chain-elongation process, core-*
565 *structure formation and secondary chain modifications by red, blue and green colours, respectively. (B) The bars represent*
566 *the diversity of the nucleotides (NT) of respective genes. The bars are colour coded as in (A). (C) Diversity of NT residues*
567 *plotted against the diversity of AA residues of respective GSL genes. Each dot is colour coded to denote genes as in (A) and*
568 *(B). The black dashed line is a linear regression line.*
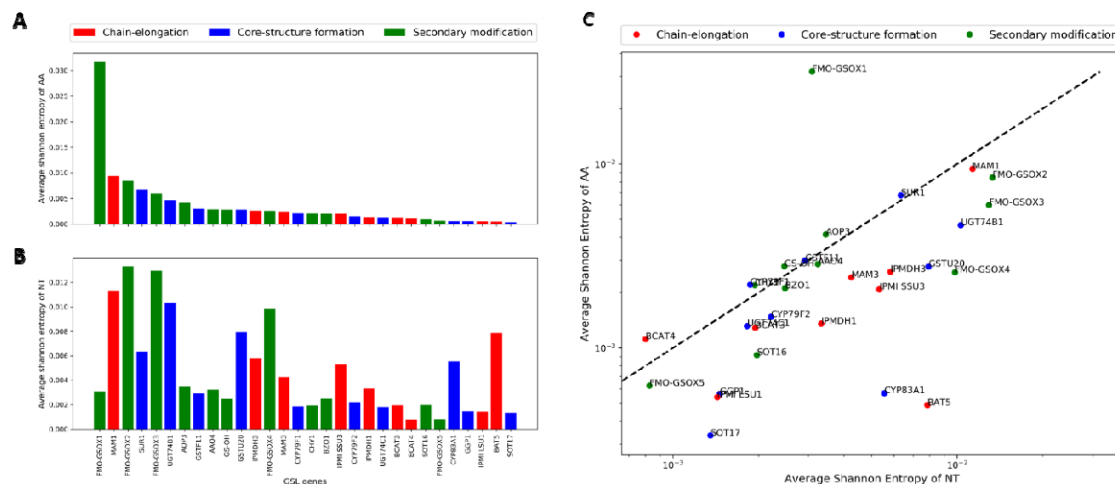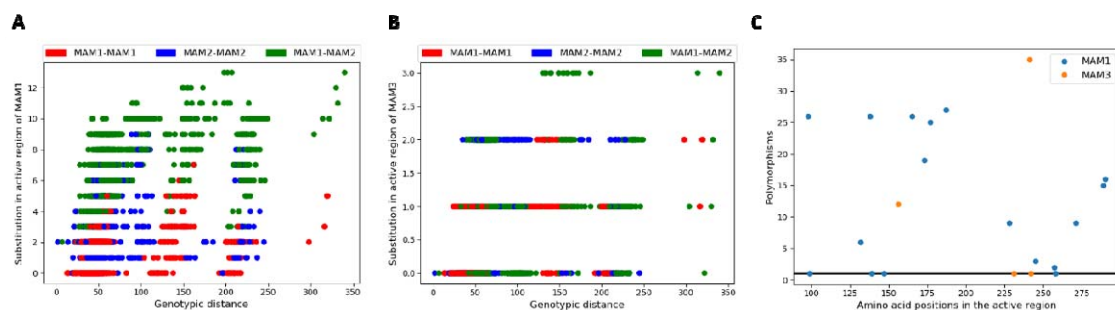


569
570

571

572

573

574 *Figure 5: Polymorphisms in the active region of MAM synthases. (A) Substitutions in the active sites of MAM1 versus the*
575 *Genotypic distances. (B) Substitutions in the active sites of MAM1 versus the Genotypic distances. (C) Polymorphisms in the*
576 *active region of MAM1 and MAM3.*



577

578

579

580  # Supplementary Figures:

581  *Figure 1: Diversity of GSL enzymes from different ecotype subgroups. The bars represent the diversity of the amino-acid*
582  *residues of respective GSL genes. The blue bars denote all 72 ecotypes, while orange and green denote ecotypes having*
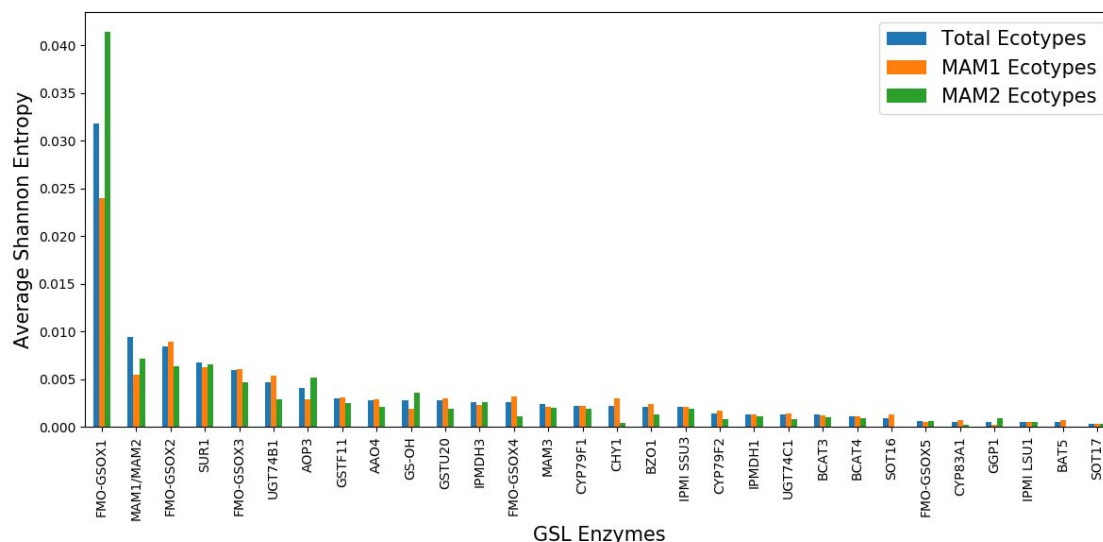583  *MAM1 and MAM2, respectively.*



584

585

586  *Figure 2: Diversity of GSL enzymes from different subgroups of ecotypes. The bars represent the diversity of nucleotide*
587  *sequences of respective GSL genes. Colour coding is same as Supplementary Figure 1.*



588